# Document made available under the Patent Cooperation Treaty (PCT)

International application number: PCT/US05/000755

International filing date: 11 January 2005 (11.01.2005)

Document type: Certified copy of priority document

Document details: Country/Office: US
Number: 60/540,634
Filing date: 30 January 2004 (30.01.2004)

Date of receipt at the International Bureau: 18 February 2005 (18.02.2005)

Remark: Priority document submitted or transmitted to the International Bureau in compliance with Rule 17.1(a) or (b)

1284004

# THE UNITED STATES OF AMERICA

## TO ALL TO WHOM THESE PRESENTS SHALL COME:

UNITED STATES DEPARTMENT OF COMMERCE

United States Patent and Trademark Office

*February 10, 2005*

**THIS IS TO CERTIFY THAT ANNEXED HERETO IS A TRUE COPY FROM THE RECORDS OF THE UNITED STATES PATENT AND TRADEMARK OFFICE OF THOSE PAPERS OF THE BELOW IDENTIFIED PATENT APPLICATION THAT MET THE REQUIREMENTS TO BE GRANTED A FILING DATE.**

**APPLICATION NUMBER: *60/540,634***
**FILING DATE: *January 30, 2004***
**RELATED PCT APPLICATION NUMBER: *PCT/US05/00755***

Certified by

Under Secretary of Commerce
for Intellectual Property
and Director of the United States
Patent and Trademark Office

PTO/SB/16 (8-00)
Approved for use through 10/31/2002. OMB 0651-0032
Patent and Trademark Office; U.S. DEPARTMENT OF COMMERCE

Under the Paperwork Reduction Act of 1995, no persons are required to respond to a collection of information unless it displays a valid OMB control number.

# PROVISIONAL APPLICATION FOR PATENT COVER SHEET

## This is a request for filing a PROVISIONAL APPLICATION FOR PATENT under 37 CFR 1.53 (c).

### INVENTOR(S)

| Given Name (first and middle [if any]) | Family Name or Surname | Residence (City and either State or Foreign Country) |
|---|---|---|
| Peng | Yin | West Windsor, New Jersey |
| Jill MacDonald | Boyce | Manalapan, New Jersey |

☐ Additional inventors are being named on the _____ separately numbered sheets attached hereto

### TITLE OF THE INVENTION (280 characters max)

**ENCODER WITH ADAPTIVE RATE CONTROL**

### CORRESPONDENCE ADDRESS

Direct all correspondence to:

☐ Customer Number

OR

Type Customer Number here

Place Customer Number Bar Code Label here

| ☒ Firm or Individual Name | Joseph S. Tripoli - Thomson Licensing Inc. |
|---|---|
| Address | PATENT OPERATIONS |
| Address | P. O. BOX 5312 |

| City | PRINCETON | State | NJ | ZIP | 08543-5312 |
|---|---|---|---|---|---|
| Country | USA | Telephone | 609-734-6834 | Fax | 609-734-6888 |

### ENCLOSED APPLICATION PARTS (check all that apply)

☒ Specification *Number of Pages*  `10`     ☐ CD(s), Number

☒ Drawing(s) *Number of Sheets*  `1`     ☐ Other (specify)

☐ Application Data Sheet. See 37 CFR 1.76

### METHOD OF PAYMENT OF FILING FEES FOR THIS PROVISIONAL APPLICATION FOR PATENT (check one)

☐ Applicant claims small entity status. See 37 CFR 1.27.

☐ A check or money order is enclosed to cover the filing fees

☒ The Commissioner is hereby authorized to charge filing fees or credit any overpayment to Deposit Account Number: `07-0832`

☐ Payment by credit card. Form PTO-2038 is attached.

FILING FEE AMOUNT ($) `160`

The invention was made by an agency of the United States Government or under a contract with an agency of the United States Government.

☒ No.

☐ Yes, the name of the U.S. Government agency and the Government contract number are: _____.

Respectfully submitted,
SIGNATURE _____

TYPED or PRINTED NAME   GUY H. ERIKSEN

TELEPHONE   (609) 734-6809

Date `01/30/04`

REGISTRATION NO. (if appropriate) `41,736`

Docket Number: `PU040032`

## USE ONLY FOR FILING A PROVISIONAL APPLICATION FOR PATENT

This collection of information is required by 37 CFR 1.51. The information is used by the public to file (and by the PTO to process) a provisional application. Confidentiality is governed by 35 U.S.C. 122 and 37 CFR 1.14. This collection is estimated to take 8 hours to complete, including gathering, preparing, and submitting the complete provisional application to the PTO. Time will vary depending upon the individual case. Any comments on the amount of time you require to complete this form and/or suggestions for reducing this burden, should be sent to the Chief Information Officer, U.S. Patent and Trademark Office, U.S. Department of Commerce, Washington, D.C., 20231. DO NOT SEND FEES OR COMPLETED FORMS TO THIS ADDRESS. SEND TO: Box Provisional Application, Assistant Commissioner for Patents, Washington, D.C. 20231.

# ENCODER WITH ADAPTIVE RATE CONTROL

Rate control is necessary in a JVT video encoder to achieve particular constant bitrates, when needed for fixed channel bandwidth applications with limited

5    buffer sizes. Avoid buffer overflow and underflow is more challenging on video content that contains sections with different complexity characteristics, for example with scene changes and dissolves.

Rate control has been studied for previous video compression standards. TMN8 [1] was proposed for H.263. The TMN8 rate control uses a frame-layer rate

10    control to select the target number of bits for the current frame and a macroblock-layer rate control to select the value of QP for the macroblocks [4].

In the frame-layer rate control, the target number of bits for the current frame is determined by

$$B = R / F - \Delta, \tag{1}$$

$$\Delta = \begin{cases} W / F, & W > Z \bullet M \\ W - Z \bullet M, & otherwise \end{cases} \tag{2}$$

$$W = \max(W_{prev} + B'-R / F, 0) \tag{3}$$

15    where $B$ is the target number of bits for a frame, $R$ is the channel rate in bits per second, $F$ is the frame rate in frames per second, $W$ is the number of bits in the encoder buffer, $M$ is the maximum buffer size, $W_{prev}$ is the previous number of bits in the buffer, $B'$ is the actual number of bits used of encoding the previous frame, and $Z=0.1$ is set by default to achieve the low delay.

20    The macroblock-layer rate control selects the value of the quantization step size for all the macroblocks in a frame, so that the sum of the macroblock bits is close to the frame target $B$. The optimal quantization step size $Q_i^*$ for macroblock $i$ in a frame can be determined by

$$Q_i^* = \sqrt{\frac{AK}{\beta_i - AN_iC}} \; \frac{\sigma_i}{\alpha_i} \; \sum_{k=1}^{N} \alpha_k \sigma_k, \tag{4}$$

25    where $K$ is the model parameter, $A$ is the number of pixels in a macroblock, $N_i$ is the number of macroblocks that remain to be encoded in the frame, $\sigma_i$ is the standard deviation of the residue in the *ith* macroblock, $\alpha_i$ is the distortion weight of the *ith* macroblock, $C$ is the overhead rate, and $\beta_i$ is the number of bits left for encoding the frame by setting $\beta_1 = B$ at the initialization stage.

The TMN8 scheme is simple and is known to be able to achieve both high quality and accurate bit rate, but is not well suited to H.264. Rate-distortion optimization (RDO) (i.e., rate-constrained motion estimation and mode decision) is a widely accepted approach in H.264 for mode decision and motion estimation, where

5     the quantization parameter (QP) (used to decide $\lambda$ in the Lagrangian optimization) needs to be decided before RDO is performed [5]. But the TMN8 model requires the statistics of prediction error signal (residue) to estimate QP, which means that motion estimation and mode decision needs to be performed before QP is made, thus resulting in a chicken and egg dilemma.

10     The methods disclosed in [2] and [3] have been proposed for H.264 rate control. Method [2] has been incorporated into the JVT JM reference software release JM7.4 [4]. To overcome the chicken and egg dilemma mentioned above, method [2] uses the residue of the collocated macroblock in the most recently coded picture with the same type to predict that of the current macroblock, and method [3]

15     employs a two-step encoding, where the QP of the previous picture ($QP_{prev}$) is first used to generate the residue, and then the QP of current macroblock is estimated based on the residue. The former approach is simple, but it lacks precision. The latter approach is more accurate, but it requires multiple encoding, thus adding much complexity.

20     In this invention, we build upon the model used in TMN8 of H.263+[4]. This model uses Lagrangian optimization to minimize distortion subject to the target bitrate constraint. To adapt the model into H.264 and to further improve the performance, we have to consider several issues. First, rate-distortion optimization (RDO) (i.e., rate-constrained motion estimation and mode decision) is a widely

25     accepted approach in H.264 for mode decision and motion estimation, where the quantization parameter (QP) (used to decide $\lambda$ in the Lagrangian optimization) needs to be decided before RDO is performed [5]. But the TMN8 model requires the statistics of prediction error signal (residue) to estimate QP, which means that motion estimation and mode decision needs to be performed before QP is made, thus

30     resulting in a chicken and egg dilemma. Second, TMN8 is targeted at low delay applications. But H.264 can be used for various applications. Therefore a new bit allocation and buffer management scheme is needed for various content. Third, TMN8 adapts the QP at macroblock level. Though a constraint is made on the QP

difference (DQUANT) between current macroblock and last coded macroblock, subjective effects of large QP variations within the same picture can be observed and has a negative subjective effect. In addition, it is known that using constant QP for the whole image may save additional bits for coding DQUANT, thus achieving

5   higher PSNR for very low bit rate. Finally, H.264 uses 4x4 integer transform and if the codec uses some thresholding techniques such as in JM reference software [4], details may be lost. Therefore, it is useful to adopt the perceptual model in the rate control to maintain the details.

10                                **Preprocessing Stage**

From equation (4), we can see that the TMN8 model requires the knowledge of standard deviation of the residue to estimate QP. However, RDO requires knowledge of the QP to perform motion estimation and mode decision thus to produce the residue. To overcome this dilemma, [2] uses the residue of the

15   collocated macroblock in the most recently coded picture with the same type to predict that of the current macroblock, and [3] employs a two-step encoding, where the QP of the previous picture ($QP_{prev}$) is first used to generate the residue, and then the QP of current macroblock is estimated based on the residue. The former approach is simple, but it lacks precision. The latter approach is more accurate, but it

20   requires multiple encoding, thus adding too much complexity.

In our approach, we adopt a different method to estimate the residue, which is simpler than the method of [3], but more accurate than the method of [2]. Experiments show that a simple preprocessing stage can give a good estimation of the residue. For $I$ picture, we only test the 3 most probable intra16x16 modes

25   (vertical, horizontal and DC mode) and MSE (Mean Square Error) of the prediction residual is used to select the best mode. Only three modes are tested in order to reduce complexity. However, in other embodiments of the current invention more of fewer modes can be tested.

30   The spatial residue is then generated using the best mode. It should be noted that we use the original pixel values for intra prediction instead of reconstructed ones, simply because the reconstructed pixels are not available. For $P$ pictures, we perform a rate-constrained motion search using only the 16x16 block type and 1

reference picture. The temporal residue is generated using the best motion vector in this mode. The average QP of the previously coded picture is used to decide $\lambda$ on rate-constrained motion search. The experiment shows that by constraining the difference of QP between previous coded picture and current picture, the $\lambda$ based on $QP_{prev}$ has minor impact on motion estimation. The side advantage of this approach is that the resulted motion vectors in the preprocessing step can be used as initial motion vectors in the motion estimation during the encoding.

### Frame-layer rate control

TMN8 is targeted to low-delay and low bit rate applications, which assume to encode only P pictures after the first I picture, hence the bit allocation model as shown in equation (1) should be re-defined to adapt to the various applications which use more frequent I pictures. The QP estimation model by equation (4) can result in large QP variation within one image, thus a frame-level QP is better first estimated to put a constraint on the variation of MB QP. In addition, for very low bit rate, due to the overhead of coding the DQUANT, it may be more efficient to use a constant picture QP. So a good rate control scheme should allow rate control at both the frame-level and the MB-level.

We first propose a new bit allocation scheme. Then we shall present a simple scheme to decide a frame-level QP.

In many applications, e.g. real-time encoders, the encoder does not know the total number of frames that need to be coded beforehand, or when scene changes will occur. Thus we adopted a GOP layer rate control to allocate target bits for each picture. The H.264 standard does not actually contain Group of Pictures, but the terminology is used here to represent the distance between I pictures. The length of the GOP is indicated by $N_{GOP}$. If $N_{GOP} \rightarrow \infty$, we set $N_{GOP} = F$, which corresponds to one second's length of frames. Notation $BG_{i,j}$ is used to indicate the remaining bits in the GOP $i$ after coding picture $j$-1, equaling to

$$BG_{i,j} = \begin{cases} \min(RG_{i-1} + R/F * N_{GOP}, R/F * N_{GOP} + M * 0.2), & j = 0 \\ BG_{i,j-1} - B_{i,j-1}, & otherwise \end{cases} \quad (5)$$

In the above equation, $RG_{i-1}$ is the number of remaining bits after GOP $i$-1 is coded, given by $RG_{i-1} = R/F * N_{coded} - B_{coded}$, where $B_{coded}$ is the used bits and $N_{coded}$ is the

number of coded pictures after GOP $i$ is finished. $B_{i,j}$ and $B'_{i,j}$ is the target bits and actual used bits for frame $j$ of GOP $i$, respectively. In equation (5), we add one constraint on the total number of bits allocated for the GOP $i$ to prevent buffer overflow when the complexity level of the content varies dramatically from one GOP to another. For example, consider a scenario where the previous GOP was of very low complexity, e.g. all black, so the buffer fullness level would go quite low. Instead of allocating all of the unused bits from the previous GOP to the current GOP, the unused bits are distributed over several following GOPs by not allowing more than $0.2M$ additional bits to an individual GOP. The target frame bit $B_{i,j}$ is then allocated according to picture type. If the *jth* picture is $P$, the target bits is $B^P_{i,j} = BG_{i,j}/(K'N' + N^P)$, where $K'$ is the bit ratio between $I$ picture and $P$ picture, which can be estimated using a sliding window approach, $N'$ is the remaining number of $I$ pictures in GOP $i$ and $N^P$ is that of $P$ pictures, otherwise, $B^I_{i,j} = K'B^P_{i,j}$. Since $P$ picture are used as the references by subsequent $P$ pictures in the same GOP, we shall allocate more target bits for $P$ pictures that are at the beginning of the GOP to ensure the later $P$ pictures can be predicted from the references of better quality and the coding quality can be improved. We use a linear weighted $P$ picture target bit allocation as follows:

$$B^P_{i,j} + = R/F * 0.2 * (N_{GOP} - 2j)/(N_{GOP} - 2) \quad (6).$$

Another constraint is added to better meet that target bits for a GOP as

$$B_{i,j} + = 0.1 * B_{diff},$$

where $B_{diff,j-1} = B_{i,j-1} - B'_{i,j-1}$, and $B_{diff,j-1} = sign(B_{diff,j-1}) \min(|B_{diff,j-1}|, R/F)$.

In our rate control, we aim at 50% buffer occupancy. To prevent the buffer overflow or underflow, the target bits need to be jointly adapted with buffer level. The buffer level $W$ is updated at the end of coding each picture by equation (3). In our approach, instead of using real buffer level to adjust the target bits, a virtual buffer level $W'$ given by $W' = \max(W, 0.4M)$ is proposed. This helps prevent the scenario that if the previously coded pictures are of very low complexity such as black scenes and consume very few bits, then the buffer level will become very low. If we use the real buffer level to adjust target frame bits as in equation (7), we may allocate too many bits, which will cause QP to decrease very quickly., But after a while, when the scene returns to normal, the low QP will easily cause the buffer to overflow. Hence we need

to either increase QP dramatically or skip the frames. This causes the temporal quality to vary significant. Then we adjust the bits by buffer control as

$$B_{i,j} = B_{i,j} * (2M - W^{\cdot})/(M + W^{\cdot}) \quad (7).$$

To guarantee a minimum level of quality, we set $B_{i,j} = \max(0.6 * R/F, B_{i,j})$. To further

avoid the buffer overflow and underflow, we set buffer safety top margin $W_T$ and bottom margin $W_B$ for $I$ picture as $W_T^I = 0.75M$, and $W_B^I = 0.25M$. As for $P$ pictures, compliant with equation (5) and to allow enough buffer for the next $I$ picture in the next GOP, we set

$W_T^P = (1 - ((0.4 - 0.2)/(N - 1) * j + 0.2)) * M$, and $W_B^P = 0.1M$. The final target bits are determined as

follows. We set $W_{VT} = W + B_{i,j}$, $W_{VB} = W_{VT} - R/F$. If $W_{VT} < W_T, B- = W_{VT} - W_T$, else if

$W_{VB} < W_B$, $B+ = W_B - W_{VB}$.

We note that if a scene change detector is employed, we shall encode the picture at the scene change to be an $I$ picture and a new GOP starts from this $I$ picture. The above scheme can still be employed.

We propose a new scheme to decide frame-level QP based on equation (4). We modify (4) as

$$\hat{Q}_i = \sqrt{\frac{AK}{B - \hat{C}}} \frac{\sigma_i}{\alpha_i} \sum_{k=1}^{N} \alpha_k \sigma_k, \quad (8)$$

where $\hat{c}$ is the overhead from last coded picture with the same type, $\sigma_i$ is estimated in the preprocessing stage as in Section 3.1. Two approaches can be used to get frame-level constant QP, denoted as $QP_f$. The first approach is to set $\alpha_i = \sigma_i$, so that all the MB QPs are equal. The second method is to use the same $\alpha_i$ as that of the MB level, as defined in the next section, then use the mean, median or mode of the histogram of the $\hat{Q}_i$ values to find the $QP_f$.

In a preferred embodiment, the second method is used to better match the MB QP. The frame-level quantization step size is decided by the mean of the $\hat{Q}_i$ values, $\hat{Q}_f = \sum_{i=1}^{N} \hat{Q}_i / N$. We note that there is a conversion between the quantization parameter QP and quantization step size $Q$ by $Q = 2^{(QP-6)/6}$. To reduce the temporal quality variation between adjacent pictures, we set $QP_f = \max(QP_f' - D_f, \min(QP_f, QP_f' + D_f))$, where $QP_f'$ is the frame QP of last coded frame, and $D_f = \begin{cases} 2 & W < 0.7M \\ 4 & otherwise \end{cases}$. Since scene

changes usually cause higher buffer levels, we take advantage of temporal masking effect and set $D_f$ to be a higher value when a scene change occurs.

## MB-layer rate control

5    A first key feature in MB-layer rate control is about the adaptive selection of weighted distortion $\alpha_i$ to get a better perceptual quality. A second key feature is to reduce the variation of the MB QPs in the same picture.

For low detail content, such as an ocean wave, a lower QP is required to keep the details. But from an RDO point of view, a higher QP is preferred, because the

10    lower detail content tends to give a higher PSNR. To keep a balance, we adopt different settings of $\alpha_i$ for $I$ and $P$ pictures, respectively. For $I$ picture, a higher distortion weight is given to the MBs with less detail, so that the detail can be better retained. Accordingly, we set

$$\alpha_i = (\sigma_i + 2\sigma_{avg})/(2\sigma_i + \sigma_{avg}), \text{ where } \sigma_{avg} = \sum_{i=1}^{N} \sigma_i / N.$$

15    For $P$ picture, a higher distortion weight is given to the MBs with more residue errors as in [4]. Accordingly,

$$\alpha_i = \begin{cases} 2B/AN(1-\sigma_i)+\sigma_i, & B/AN < 0.5 \\ 1, & otherwise \end{cases}$$

In this way, better perceptual quality is maintained for $I$ picture and can be propagated to the following $P$ pictures, while higher objective quality is still kept as in

20    [4]. To prevent large variation of the quality inside one picture, we set $QP_i = \max(QP_f - 2, \min(QP_i, QP_f + 2))$. If a frame level rate control is used, $QP_i = QP_f$.

## Virtual frame skipping

After encoding one picture, we shall update $W$ by equation (3). If $w > 0.9M$, the

25    next frame is virtually skipped until the buffer level is below $0.9M$. Virtual frame skipping is to code every MB in the $P$ picture to be SKIP mode. In this way, we can syntactically keep a constant frame rate. If the current frame is decided to be a virtual skipped frame, we set $QP_f = QP_f' + 2$.

In summary, our rate control scheme consists of the following steps:

30    preprocessing, frame target bits allocation and frame-level constant QP estimation,

MB-level QP estimation, buffer updates and virtual frame skipping control. Our approach can allow both frame-level and MB-level rate control.

### Advantages/Benefits From Various Embodiments

5      **1.** The use of mean/median/mode of initial macroblock QP estimates to select frame level QP.

      **2.** When the selected frame level QP is used in the calculation of the individual macroblock QPs.

      **3.** When performing intra prediction using a subset of the allowable intra-
10      prediction modes to form the residue that is used in the QP selection process.

      **4.** Use of a small number of intra-prediction modes (three (3), for example).

      **5.** When a previous GOP was coded with a large number of unused bits, limiting the additional bits allocated to the current GOP to a predetermined threshold.

      **6.** When a virtual buffer level instead of real buffer level is to used for buffer
15      control.

20     Figure 1 is a video encoder and is indicated generally by the reference numeral 100. An input to the encoder 100 is connected in signal communication with a non-inverting input of a summing junction 110. The output of the summing junction 110 is connected in signal communication with a block transform function 120. The transformer 120 is connected in signal communication with a quantizer 130. The output of the quantizer 130 is connected in signal communication with a variable length coder ("VLC") 140, where the output of the VLC 140 is an externally available output of the encoder 100.
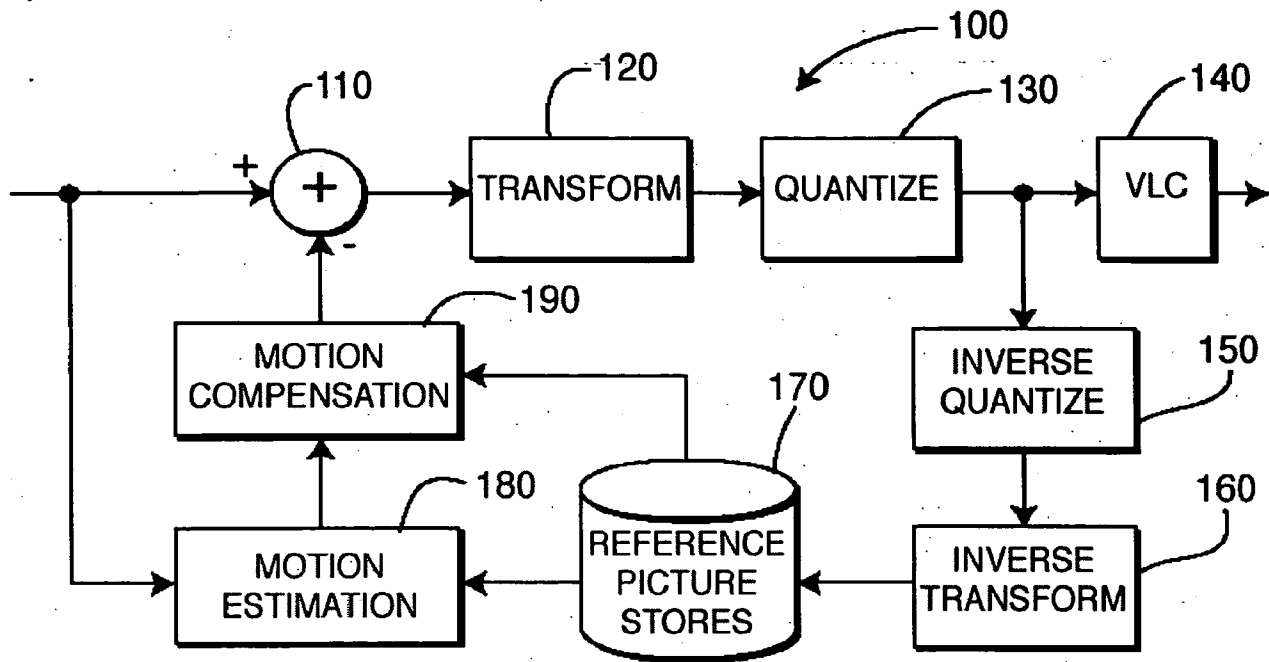
25     The output of the quantizer 130 is further connected in signal communication with an inverse quantizer 150. The inverse quantizer 150 is connected in signal communication with an inverse block transformer 160, which, in turn, is connected in signal communication with a reference picture store 170. A first output of the reference picture store 170 is connected in signal communication with a first input of
30     a motion estimator 180. The input to the encoder 100 is further connected in signal communication with a second input of the motion estimator 180. The output of the

motion estimator 180 is connected in signal communication with a first input of a motion compensator 190. A second output of the reference picture store 170 is connected in signal communication with a second input of the motion compensator 190. The output of the motion compensator 190 is connected in signal

5    communication with an inverting input of the summing junction 110.

# References

[1] J. Ribas-Corbera and S. Lei, "Rate Control in DCT Video Coding for Low-Delay Communications", IEEE Trans. *CSVT*, , Feb., 1999.

[2] S. Ma and W. Gao, "Adaptive Rate Control with HRD Consideration," *JVT-H014*, 8th meeting, Geneva, May,2003

[3] Z. He and T. Chen, "Linear Rate Control for JVT Video Coding," *ITRE*, 2003

[4] K. Suehring, JVT JM reference software, http://bs.hhi.de/~suehring/tml/download/

[5] T. Wiegand, G. J. Sullivan, G. Bjontegaard, and A. Luthra, "Overview of the H.264/AVC Video Coding Standard," IEEE. Trans. *CSVT*, pp. 560-576, July, 2003.

[6] P. Yin, H. C. Tourapis, A. Tourapis and jJ. Boyce, "Fast mode decision and motion estimation for JVT/H.264", *ICIP2003*

**FIG. 1**